

A proposal of a criterion for collision resistance of hash functions

Dai Watanabe¹ and Hirotaka Yoshida¹

Systems Development Laboratory, Hitachi, Ltd.

{daidai,hyoshida}@sdl.hitachi.co.jp

Abstract

In this paper we revisit the techniques for collision attacks and study the relation between maximum differential characteristic probability and a limit of applicability of collision attack. We show that a cryptographic hash function is secure against collision attacks using a single message block based on differential attack if the inequality $p_D < (1 - e^{-1})2^{-n_m-1}$ is satisfied, where n_m is an input length of a compression function and p_D is the maximum differential characteristic probability.

Keywords. Hash function, Collision attack, Differential characteristic

1 Introduction

A hash function is a cryptographic primitive which compresses data of arbitrary length into a fixed length bit strings. A hash function play a crucial role especially in authentication mechanisms such as a digital signature and a message authentication code so that it is required to be highly secure. The basic security requirement for a hash function is so called *collision resistance* which is the difficulty to find two distinct inputs whose outcomes are the same.

For long time it has been unclear what collision resistance is, and how to evaluate the strength against collision attacks. Addition to few examples of the algorithms and their evaluations, the essential fact that there is no secret (a key) in hash calculations confused many researchers. At last, it was not

clear the advantage of the fact that an attacker can know all intermediate values in calculating an output. This fact is the most different assumption for an attacker from block cipher's case.

However Wang *et al.* showed in the last two years that almost all the currently proposed hash functions (including widely used MD5 and SHA-1) is weak against their collision attacks [16, 17, 18, 19]. Additionally Biham *et al.* provided a technique to improve the complexity of collision attacks and applied it to SHA-0 and SHA-1 [1, 2]. Both of their attacks are an application of differential attack proposed by Biham and Shamir which was originally applied to the block cipher DES for recovering a secret key [3]. With helps of these newer proposed techniques and their applications, the standing position of the probabilistic approach in collision attack begun to be clear.

In this paper, we revisit the known techniques for collision attacks and try to clarify the relationship between collision attacks and naive differential attacks. As a result we propose a criterion of collision resistance from a viewpoint of differential probability.

The organization of this paper is as follows: Firstly we introduce the basic terminologies in Sect. 2. Secondly Dobbertin, Biham, and Wang's collision attacks are revisited in Sect. 3. In Sect. 4, we observe the relation between naive differential attacks and collision attacks and propose a criterion of collision resistance. In Sect. 5 the adequacy of the proposed criterion is discussed. Finally we conclude the discussion in Sect. 6

2 Preliminary

In this section we give a brief explanation of terminologies used in this paper.

2.1 How to construct a hash function

A general method to construct a hash function which deals with a message of arbitrary length is to divide a message into several blocks of fixed length and to process them sequentially. A function h which processes a message block of fixed length is called a *compression function*. The most widely used method to process is so-called *Merkle-Damgård strengthening*, which is defined as follows:

$$H = H_n, H_i = h(H_{i-1}, M_i),$$

where a message M is divided into n blocks M_1, \dots, M_n . Merkle and Damgård independently proved that this chaining construction is secure as a hash function if the underlying compression function is secure [5, 12].

2.2 Security requirements for a hash function

Following three conditions are the security requirements for a hash function.

One-Wayness For any hash value y it is difficult to find an input x such that $Hash(x) = y$.

Second Pre-image Resistance For any input x it is difficult to find a distinct input x' such that $Hash(x) = Hash(x')$

Collision Resistance It is difficult to find a pair of inputs $(x, x'), x \neq x'$ such that $Hash(x) = Hash(x')$.

Throughout this paper we deal with only the third condition.

The security requirements for a compression function are almost the same as that for a hash function except a point that the input has a context. In detail, the input of a compression function as a sub function of a hash function is divided into two parts,

an intermediate hash value H_{i-1} which is the output of the previous application of the compression function and a message block M_i . If the underlying compression function is an ideal function, the intermediate hash value H_{i-1} is randomly distributed. Hence the attacker is usually assumed to be able to control only message inputs M_i when considering the security of a compression function as a sub function of a hash function. On the other hand, when the security of a compression function itself, the attacker can control not only message input M_i but also a hash input H_{i-1} . The collision resistance under this scenario is called *Pseudo-Collision Resistance*.

Basically the security of a hash function depends on the length of the output, called the hash length. Let the hash length be n_h bits, then it is necessary to calculate the target hash function about 2^{n_h} times to find a pre-image or a second pre-image by brute force. Only for the collision resistance a generic attack which is much faster than brute force is known. The attack is called the birthday attack because it is based on the famous *birthday paradox* which clarify a significant property of a random set. The birthday attack shows that it is possible to find a collision with about $2^{n_h/2}$ inputs. This fact claims that the hash length should be twice as large as that of a block length (of a block cipher) used in the same system. For more detail of a birthday attack, please refer to [11] for example.

2.3 Differential attack

Differential attack was proposed by Biham and Shamir for the attack on the block cipher DES (Data Encryption Standard) [3]. In this subsection we give a brief description of differential attack.

Let X, Y be groups and '+' be the operation on them (for example arithmetic addition or xoring). For a map f from X to Y , the *differential* of f by the difference Δ_x is defined as follows:

$$\Delta f(x, \Delta_x) := f(x + \Delta_x) - f(x).$$

If f is an ideal random function, the function Δf must be random independent of the input difference

Δ . The basic idea of differential attack is to study the distribution of Δf depending on the input difference Δ_x to distinguish f from a truly random function.

From now on let X, Y be vector space on $\text{GF}(2)$ of dimensions n_x, n_y , respectively. Let f be a map from X to Y . A *differential probability* associated with the input difference Δ_x and the output difference Δ_y is defined as follows:

$$DP(f)(\Delta_x, \Delta_y) := \frac{|\{x \in X | f(x + \Delta_x) - f(x) = \Delta_y\}|}{2^{n_x}}.$$

The *maximum differential probability* is the maximum value of the differential probability with all pairs of non-zero input and output differences and defined as follows:

$$DP_{\max}(f) := \max_{\Delta_x \neq 0, \Delta_y} DP(f)(\Delta_x, \Delta_y).$$

If the function f is an ideal random function, $DP_{\max}(f) \approx 2^{-n_y}$ is satisfied.

It is difficult in practice to calculate the maximum differential probability of the real block ciphers or hash functions because their input and output bit lengths are too large. These functions are usually designed in cascading style, i.e., they can be usually decomposed to sub functions f_i and the output of the sub function f_1 is input to the next function f_2 , and so on. In such case the *maximum differential characteristic probability* which is defined by the multiplying the maximum differential probabilities of f_i for all i . It is often applied to evaluate the lower bound of the maximum differential probability.

Let f be a cascading function such that $f = f_r \circ f_{r-1} \circ \dots \circ f_1$, then the maximum differential characteristic probability is defined as follows:

$$DCP_{\max}(f) := \max_{\Delta_x \neq 0, \Delta_y} \prod_{\substack{0 < i \leq r, \\ \Delta_x = \Delta_0, \Delta_r = \Delta_y}} DP(f)(\Delta_{i-1}, \Delta_i)$$

The sequence of differences $(\Delta_0, \Delta_1, \dots, \Delta_r)$ which gives the maximum differential characteristic probability is called the best differential path of the function f .

3 Known collision attacks

All known collision attacks are the application of differential attack. In these attacks firstly the differential path whose output difference is equal to zero is fixed. Let the differential characteristic probability of the path be p . Then it is expected that a colliding pair is found if about p^{-1} trials are executed. Hence if there is a differential path with probability p satisfying $1/2 \cdot p^{-1} < 2^{n_h/2}$, the differential attack can effectively find a collision compared with birthday attack. In other words, the collision resistance of the target hash function (or the compression function) is not sufficient. This is the basic idea of differential based collision attacks.

The important known collision attacks on certain hash functions are presented by Dobbertin, Biham, and Wang, and all of them are applications of basic differential attack described above. In this section their attacks are revisited.

3.1 Dobbertin's technique

Dobbertin blazes a way on a collision attack by studying the early proposed hash function such as MD4, MD5, and RIPEMD [8]. The outline of Dobbertin's collision attack is described as follows:

Algorithm 1 Dobbertin's collision attack

Step 1. Fix a differential path whose output difference is zero.

Step 2. For the first several steps write up the equations by using intermediate variables to make the behavior of differences deterministic.

Step 3. Solve the system of equations and execute random testing using the solutions.

Dobbertin's attack narrows down the input set satisfying the differential path for several steps by solving the system of equations (consisting of 32-bit-wise logical operations and arithmetic addition). As a result the differential characteristic probability of the given path on the set used in Step 3 is larger than the random testing so that the complexity of the collision attack is reduced.

3.2 Biham’s technique

Biham and Chen defined a concept of *neutral bit* and reduced the calculation complexity of collision attack [1, 2]. The outline of Biham’s collision attack is described as follows:

Algorithm 2 Biham’s collision attack

Step 1. Fix a differential path whose output difference is zero.

Step 2. Find an input which satisfies the differential path for the first several steps by random testing. Denote the input by P_0 .

Step 3. Let e_j be a vector whose bits are zero except in the j -th bit position and \mathcal{N} be the set consisting of vectors e_j which does not have influence on the differences for the first several steps. The elements of \mathcal{N} are called *neutral bits*.

Step 4. Execute random testing to find a collision by choosing the inputs from the set $\{P_0 + \varepsilon | \varepsilon = \sum e_j, e_j \in \mathcal{N}\}$.

This attack provides a generic methods to gather a set whose elements are satisfying the differential path for several steps. As a result the differential characteristic probability of the given path on the set used in Step 3 is larger than the one in the random testing so that the calculation complexity of the collision attack is reduced.

3.3 Wang’s technique

Wang *et al.* proposed a technique called *message modification* to reduce the calculation complexity and applied them to currently widely used hash algorithms such as MD5 and SHA-1 [16, 17, 18, 19]. The outline of Wang’s collision attack is described as follows:

The basic idea of Wang’s technique is almost the same as Dobbertin’s technique. However it does not solve the system of equations of sufficient conditions. Instead, it modified the input in the online manner. Additionally the attack chooses the better differential path than what was used by Dobbertin and studies their bitwise sufficient conditions. Because

Algorithm 3 Wang’s collision attack

Step 1. Fix a differential path whose output difference is zero.

Step 2. For each step operation whose output difference is probabilistic write up the conditions by using intermediate variables to make the behavior of differences deterministic.

Step 3. Choose inputs randomly and modify some bits according to the conditions written up in Step 2. Continue Step 3 until a collision pair is found.

of these improvements, Wang’s collision attacks are much more efficient than Dobbertin’s attacks.

4 A rough criterion of collision resistance

The observation in the previous section clarifies that the basic strategy of known collision attacks based on differential attack is all the same, which is to find an input sub space with the elements satisfying a certain differential path is satisfied with higher probability than randomly chosen input. In this section we give a simple relational expression between differential probability and collision resistance based on those observations and sum up it as a criterion of collision resistance.

4.1 Collision resistance and differential probability

All known collision attacks consist of two phases. Firstly they search for the (almost) best collision-producing differential path. Next they search for the adequate inputs which satisfies a part of the differential path and try to find a concrete collision pair with the input set. The latter process dramatically improves the required number of trials for the attack compared to what is expected from the differential characteristic probability. This is the essential difference between differential attack and collision attack. The question is how to estimate the efficiency

of the latter process. In this section we deal with this problem.

Let the input (message) length and the output (hash) length of the compression function h be n_m bits and n_h bits respectively. The differential probability in the definition is the ratio of the inputs with the input difference whose corresponding output difference are expected value. Hence a collision-producing differential path with probability p means that about $2^{n_m} \cdot p$ of inputs are expected to satisfy the path (so it collides). As a result, if $p < 2^{-n_m}$ is satisfied it looks difficult to find a collision with the input differential because the expected value of collision-producing pair is less than 1.

Meanwhile the efficiency of the attack is usually represented by the success probability of the attack with a number of trial q as a parameter. We are going to follow this manner and compare birthday attack with differential based collision attack regarding efficiency.

From now on we assume that the target compression function can be decomposed into r sub functions h_i , i.e., $h = h_r \circ h_{r-1} \circ \dots \circ h_1$. In the attack we fix the best differential path of the compression function h with its differential characteristic probability p_D , where the differential characteristic probability of each sub function h_i is given by p_{D_i} . Let U_i be the set of the inputs satisfying the fixed differential path on the sub function h_i . U_i includes about $2^{n_m} \cdot p_{D_i}$ elements and the differential path holds for the whole compression function h with high probability

$$p_D|_{U_i} = \prod_{j \neq i} p_{D_j} = p_D \cdot \frac{2^{n_m}}{\#U_i}.$$

By generalizing the discussion above and choosing the input sub space U adequately, the differential characteristic probability on the sub space U can be expressed as follows:

$$p_D|_U = p_D \cdot \frac{2^{n_m}}{\#U}.$$

The probability to find a collision with q elements in U can be approximated by $p_D|_U \cdot q$ if q is sufficiently smaller than $p_D|_U^{-1}$.

On the other hand the probability to find a collision with q inputs is generally estimated by birthday paradox, and is approximately $1 - \exp(-q^2/2^{n_h+1})$. Therefore the differential based collision attack is more effective than birthday attack iff the following inequality is satisfied:

$$1 - e^{-\frac{q^2}{2^{n_h+1}}} < p_D|_U \cdot q \quad (1)$$

The number of the elements in the trial space $\#U$ is necessary not to be smaller than q , so that $\#U \geq q$. With this the inequality (1) can be transformed as follows:

$$\begin{aligned} p_D &> 2^{-n_m} (1 - e^{-\frac{q^2}{2^{n_h+1}}}) \\ &\geq 2^{-n_m} (1 - e^{-1}) q^2 / 2^{n_h+1} \\ &= (1 - e^{-1}) \cdot 2^{-n_m-n_h-1} \cdot q^2. \end{aligned}$$

By evaluating the maximum value of the left part of the inequality the discussion is summarized as follows:

Theorem 1 (A criterion of collision resistance) *Let n_m be the message input length of the compression function h and p_D be the maximum differential characteristic probability of h . Then h is secure against differential based collision attack using a single message block if $p_D < (1 - e^{-1})2^{-n_m-1}$ is satisfied.*

Theorem 1 means that the collision resistance can be represented by the maximum differential characteristic probability. It is an interesting point that the theorem indicates the collision resistance depends on message input length rather than hash length.

In the discussion above we assumed that the hash input of n_h bits is fixed and the attacker can control only message input of n_m bits. However it is well known this condition is relaxant if the target hash function adopts Merkle-Damgård strengthening. In this case the hash function is cascading compression functions so that the attacker can get additional space of q_1 elements for the target compression function by calculating the outputs of the

previous compression function. This step can be executed independently of the collision search for the target compression function. Let q_1 be the number of trials in the first step and q_2 be the number of trials in the second step (discussed in Theorem 1). Then the following inequality is the necessary and sufficient condition that the differential attack works more efficiently than birthday attack.

$$1 - e^{-\frac{(q_1+q_2)^2}{2^{n_h+1}}} < p_D |U \cdot q_1 \cdot q_2.$$

This inequality is transformed in the same manner with the above and the result is summarized to Theorem 2.

Theorem 2 (A Criterion of collision resistance for a hash function using MD-strengthening) *Let n_m be the message input length of the compression function h and p_D be the maximum differential characteristic probability of h . Then the hash function based on h and MD-strengthening is secure against differential based collision attack using multi message blocks if $p_D < (1 - e^{-1})2^{-n_m - n_h/2 - 1}$ is satisfied.*

4.2 Pseudo-collision resistance and differential probability

Pseudo-collision resistance against differential based collision attack can be discussed in the same manner. In pseudo-collision attack the attacker can choose any input bits so that he can control $n_m + n_h$ bits of input.

Theorem 3 (A criterion of pseudo-collision resistance) *Let n_m be the message input length of the compression function h and p_D be the maximum differential characteristic probability of h . Then h is secure against differential based pseudo-collision attack if $p_D < (1 - e^{-1})2^{-n_m - n_h - 1}$ is satisfied.*

5 Ambiguity of the proposed criterion

In this section the accuracy and other problems of the proposed criterion are clarified.

5.1 Accuracy of the criterion

The assumption in the discussion in the previous section is that the attacker can find a collision if there is a collision. But this assumption is not plausible in real cases. There are the big difference between what the proposed criterion claims and what the known collision attacks show.

Table 1 shows the calculation complexities of Wang's collision attacks and their differential characteristic probability of the differential path used in the attack (the differential characteristic probabilities are estimated by counting up their sufficient conditions). Holding SHA-1 up as an example, the message block length is 512 bits and there is a differential path whose differential characteristic probability is 2^{-247} . Ideally it is possible to choose the input set on which the differential path is satisfied with probability 1, however the attack presented in [19] provides an input set on which pairs of input collide with probability 2^{-68} .

5.2 Collision attack and Markov assumption

In the definition of differential characteristic probability the target function is assumed to be a Markov cipher, i.e., the probabilistic events on sub functions are independent each other. This assumption is valid if the target function is a block cipher. In the evaluation of the encryption function of a block cipher, its key scheduling function is usually ignored and all sub keys are assumed to be random. This manner is originated in the common understanding that the attacker cannot know the information of a key bit.

However in the case of a hash function, the attacker can control all input. Additionally collision attacks are the application of differential attack and

Table 1: The message block length of the compression functions and their differential characteristic probabilities

Algorithm	Message block length (bit)	Differential characteristic probability	Complexity of collision attack	Reference
MD4	512	2^{-122}	2^{-2}	[16]
MD5	512	2^{-258}	2^{-39}	[17]
RIPEMD	512	2^{-124}	2^{-18}	[16]
SHA-0	512	2^{-218}	2^{-39}	[18]
SHA-1	512	2^{-247}	2^{-68}	[19]

their applications are mainly searching for the adequate input set, whose elements satisfy the differential path with much higher probability than random testing. Under this condition the target function no longer holds a Markov property. For example, as a result of flipping some bits of the input in Wang’s technique, some input of some sub functions changes some of their bits. Summary of these fact indicates that the Theorems claimed in the previous section are not always satisfied.

5.3 Few more problems

Now we discuss the way the standard hash functions should be. What claimed throughout this paper is their underlying compression functions should satisfy at least the condition described in Theorem 1 and hopefully the condition described in Theorem 3. So far there are some brief reports on evaluation of the differential probabilities of SHA-256, -384, -512, which are the new hash standards established by NIST (National Institute of Standard and Technologies) [10]. Their evaluations show some upper bounds of differential characteristic probability, but they does not look tight. So the first thing we should do is to give more detailed upper (or lower) bounds of differential probability for SHA2-family. For new proposals from now it is desirable to satisfy at least the condition claimed in Theorem 1.

On the other hand we acknowledge the criteria proposed in this paper is not perfect. As discussed in the section 5.1, it is usually difficult to satisfy

all sufficient conditions in Wang’s technique. This difficulty looks to show the difference of the strength against collision attack between the cases of MD5 and SHA-1. If a new criterion which clarifies the difficulty to narrow down the input set is defined, it is rather not preferable to reduce the complexity of collision attacks to differential probability as in this paper, and more flexible designs will be allowed.

Additionally it can be not necessary to fix intermediate differences in a differential path to discuss collision attacks. For example Dobbertin reported that the experimental result shows his collision attack on MD4 works better than expected [6]. More precisely, the success probability of finding collision is higher than what is expected from the differential characteristic probability. He analyzed that this deviance arises from the lack to evaluate other possible differential paths. This example indicates that it is not necessary to fix intermediate differences in a differential path for random testing. Biham’s technique also maintain this relaxation. Therefore the criteria proposed in this paper may not be correct if there is a big difference between the maximum differential characteristic probability and the maximum differential probability.

6 Conclusion

In this paper the techniques for collision attacks are revisited and the relation between maximum differential characteristic probability and a limit of applicability of collision attack are clarified. As a result

we showed that a cryptographic hash function is secure against collision attacks based on differential attack if the inequality $(1 - e^{-1})p_D < 2^{-nm-1}$ is satisfied. The study in this paper ignores some certain conditions for the simple discussion so that the resultant criterion should be dealt with care. However we wish it will be a help to understand collision resistance of a hash function.

Acknowledgement

This work was supported by a consignment research from the National Institute on Information and Communications Technology (NICT), Japan.

References

- [1] E. Biham and R. Chen, "Near collision for SHA-0," *Advances in Cryptology, CRYPTO 2004*, LNCS 3152, pp. 290–305, Springer-Verlag, 2004.
- [2] E. Biham, R. Chen, and A. Joux, "Collisions of SHA-0 and reduced SHA-1," *Advances in Cryptology, Eurocrypt 2005*, LNCS 3494, pp. 37–57, Springer-Verlag, 2005.
- [3] E. Biham and A. Shamir, *Differential Cryptanalysis of the Data Encryption Standard*, Springer-Verlag, 1993.
- [4] F. Chabaud and A. Joux, "Differential collisions in SHA-0," *Advances in Cryptology, CRYPTO'98*, LNCS 1462, pp. 56–71, Springer-Verlag, 1999.
- [5] I. Damgård, "A Design principle for hash functions," *Advances in Cryptology, CRYPTO'89*, LNCS 435, pp. 416–427, Springer-Verlag, 1990.
- [6] H. Dobbertin, "Cryptanalysis of MD4," *Journal of CRYPTOLOGY*, No. 11, pp. 253–271, 1998.
- [7] National Institute of Standards and Technologies, *Secure Hash Standard*, Federal Information Processing Standards Publication, FIPS 180, 1993.
- [8] H. Dobbertin, "Cryptanalysis of MD4," *Fast Software Encryption, FSE'96*, LNCS 1039, pp. 53–69, Springer-Verlag, 1996.
- [9] National Institute of Standards and Technologies, *Secure Hash Standard*, Federal Information Processing Standards Publication, FIPS 180-1, 1995.
- [10] National Institute of Standards and Technologies, *Secure Hash Standard*, Federal Information Processing Standards Publication, FIPS 180-2, 2002.
- [11] A. Menezes, P. van Oorschot, S. Vanstone, *Handbook of Applied Cryptography*, CRC Press, 1996.
- [12] R. Merkle, "One-way hash functions and DES," *Advances in Cryptology, CRYPTO'89*, LNCS 435, pp. 347–361, Springer-Verlag, 1990.
- [13] R. Rivest, *The MD4 message-digest algorithm*, RFC 1186, 1990.
- [14] R. Rivest, *The MD5 message-digest algorithm*, RFC 1321, 1992.
- [15] RIPE, *Integrity Primitives for Secure Information Systems. Final Report of RACE Integrity Primitives Evaluation (RIPE-RACE 1040)*, LNCS 1007, Springer-Verlag, 1995.
- [16] X. Wang, X. Lai, D. Feng, H. Chen, and X. Yu, "Cryptanalysis of the hash functions MD4 and RIPEMD," *Advances in Cryptology, Eurocrypt 2005*, LNCS 3494, pp. 1–18, Springer-Verlag, 2005.
- [17] X. Wang and H. Yu, "How to break MD5 and other hash functions," *Advances in Cryptology, Eurocrypt 2005*, LNCS 3494, pp. 19–36, Springer-Verlag, 2005.
- [18] X. Wang, H. Yu, and L. Yin, "Efficient collision search attacks on SHA-0," *Advances in Cryptology, CRYPTO 2005*, LNCS 3621, pp. 1–16, Springer-Verlag, 2005.

- [19] X. Wang, L. Yin, and H. Yu, "Finding collision in the full SHA-1," *Advances in Cryptology, CRYPTO 2005*, LNCS 3621, pp. 17-37, Springer-Verlag, 2005.